

8-2008

Choquet Integrals and OWA Criteria as a Natural (and Optimal) Next Step After Linear Aggregation: A New General Justification

Francois Modave

Martine Ceberio

The University of Texas at El Paso, mceberio@utep.edu

Vladik Kreinovich

The University of Texas at El Paso, vladik@utep.edu

Follow this and additional works at: https://scholarworks.utep.edu/cs_techrep



Part of the [Computer Engineering Commons](#)

Comments:

Technical Report: UTEP-CS-08-28a

To appear in Alexander Gelbukh and Eduardo F. Morales (eds.), *Proceedings of the 7th Mexican International Conference on Artificial Intelligence MICAI'08*, Mexico City, Mexico, October 27-31, 2008, Springer Lecture Notes on Artificial Intelligence, 2008, Vol. 5317, pp. 741-753.

Recommended Citation

Modave, Francois; Ceberio, Martine; and Kreinovich, Vladik, "Choquet Integrals and OWA Criteria as a Natural (and Optimal) Next Step After Linear Aggregation: A New General Justification" (2008). *Departmental Technical Reports (CS)*. 109.
https://scholarworks.utep.edu/cs_techrep/109

This Article is brought to you for free and open access by the Computer Science at ScholarWorks@UTEP. It has been accepted for inclusion in Departmental Technical Reports (CS) by an authorized administrator of ScholarWorks@UTEP. For more information, please contact lweber@utep.edu.

Choquet Integrals and OWA Criteria as a Natural (and Optimal) Next Step After Linear Aggregation: A New General Justification

François Modave, Martine Ceberio, and Vladik Kreinovich

Department of Computer Science
University of Texas at El Paso, El Paso, TX 79968, USA,
{fmodave,mceberio,vladik}@utep.edu

Abstract. In multi-criteria decision making, it is necessary to aggregate (combine) utility values corresponding to several criteria (parameters). The simplest way to combine these values is to use linear aggregation. In many practical situations, however, linear aggregation does not fully adequately describe the actual decision making process, so non-linear aggregation is needed.

From the purely mathematical viewpoint, the next natural step after linear functions is the use of quadratic functions. However, in decision making, a different type of non-linearities are usually more adequate than quadratic ones: non-linearities like OWA or Choquet integral that use min and max in addition to linear combinations. In this paper, we explain the empirically observed advantage of such aggregation operations.

1 Introduction

One of the main purposes of Artificial Intelligence in general is to incorporate a large part of human intelligent reasoning and decision-making into a computer-based systems, so that the resulting intelligent computer-based systems help users in making rational decisions. In particular, to help a user make a decision among a large number of alternatives, an intelligent decision-making systems should select a small number of these alternatives – alternatives which are of the most potential interest to the user.

For example, with so many possible houses on the market, it is not realistically possible to have a potential buyer inspect all the house sold in a given city. Instead, a good realtor tries to find out the buyer's preferences and only show him or her houses that more or less fit these preferences. It would be great to have an automated system for making similar pre-selections.

To be able to make this selection, we must elicit the information about the user preferences.

In principle, we can get a full picture of the user preferences by asking the user to compare and/or rank all possible alternatives. Such a complete description of user preferences may be sometimes useful, but in decision making applications, such an extensive question-asking defeats the whole purpose of intelligent

decision-making systems – to avoid requiring that the the user make a large number of comparisons.

The existing approach to this problem is called *multi-criteria decision making* (MCDM). The main idea behind this approach is that each alternative is characterized by the values of different parameters. For example, the buyer’s selection of a house depends on the house’s size, on its age, on its geographical location, on the number of bedrooms and bathrooms, etc. The idea is to elicit preferences corresponding to each of these parameters, and then to combine these single-parameter preferences into a reasonable model for describing the user’s choice.

In the standard decision making theory, preferences are characterized by assigning, to each alternative, a numerical value called its *utility*. In these terms, the multi-criteria decision making approach means that we try to combine single-variable utility values $u_1(x_1), \dots, u_n(x_n)$ characterizing the user’s preferences over individual parameters x_1, \dots, x_n into a utility value $u(x_1, \dots, x_n)$ that characterizes the utility of an alternative described by the values (x_1, \dots, x_n) .

In the first approximation, it makes sense simply to add the individual utility values with appropriate weights, i.e., to consider linear aggregation

$$u(x_1, \dots, x_n) = w_1 \cdot u_1(x_1) + \dots + w_n \cdot u_n(x_n).$$

In many practical situations, linear aggregation works well, but in some cases, it leads to counterintuitive conclusions. For example, when selecting a house, a user can assign certain weights to all the parameters characterizing different houses, but the user may also has absolute limitations: e.g., a user with kids may want a house with at least two bedrooms, and no advantages in location and price would entice her to buy a one-bedroom house. To describe such reasonable preferences, we must therefore go beyond linear aggregation functions.

From the purely mathematical viewpoint, the inadequacy of a linear model is a particular example of a very typical situation. Often, when we describe the actual dependence between the quantities in physics, chemistry, engineering, etc., a linear expressions $y = c_0 + c_1 \cdot x_1 + \dots + c_n \cdot x_n$ is a very good first approximation (at least locally), but to get a more accurate approximations, we must take non-linearity into account. In mathematical applications to physics, engineering, etc., there is a standard way to take non-linearity into account: if a linear approximation is not accurate enough, a natural idea is to use a quadratic approximation $y \approx a_0 + \sum_{i=1}^n c_i \cdot x_i + \sum_{i=1}^n \sum_{j=1}^n c_{ij} \cdot x_i \cdot x_j$; if the quadratic approximation is not sufficient accurate, we can use a cubic approximation, etc.; see, e.g., [3].

At first glance, it seems reasonable to apply a similar idea to multi-criteria decision making and consider quadratic aggregation functions

$$u \stackrel{\text{def}}{=} u(x_1, \dots, x_n) = u_0 + \sum_{i=1}^n w_i \cdot u_i(x_i) + \sum_{i=1}^n \sum_{j=1}^n w_{ij} \cdot u_i(x_i) \cdot u_j(x_j).$$

Surprisingly, in contrast to physics and engineering applications, quadratic approximation do not work as well as approximations based on the use of piece-wise

linear functions, such as the OWA operation $u = w_1 \cdot u_{(1)} + \dots + w_n \cdot u_{(n)}$, where $u_{(1)} = \max(u_1(x_1), \dots, u_n(x_n))$ is the largest of n utility values $u_i(x_i)$, $u_{(2)}$ is the second largest, \dots , and $u_{(n)} = \min(u_1(x_1), \dots, u_n(x_n))$ is the smallest of n utility values; see, e.g., [9].

In our own research, we have applied OWA and we have also applied similar piece-wise linear operations (based on the so-called Choquet integral [4]), and we also got good results – better than quadratic approximations; see, e.g., [1] and references therein. Similar results have been obtained by others. For quite some time, why piece-wise approximations are better than quadratic ones remains a mystery to us – and to many other researchers whom we asked this question. Now, we finally have an answer to this question – and this answer is presented in the current paper.

Thus, the paper provides a new justification of the use of piece-wise aggregation operations in multi-criteria decision making – a justification that explains why these aggregation operations are better than the (seemingly more natural) quadratic ones.

The structure of this paper is as follows. To explain our answer to the long-standing puzzle, we need to recall the properties of the utility functions. The needed properties of utility functions are described in Section 2. Readers who are already well familiar with the standard decision making theory (and with the corresponding properties of utility functions) can skip this section and proceed directly to Section 3. In Section 3, we explain why quadratic aggregation operations are less adequate than OWA and Choquet operations: because the basic properties of utility functions lead to the need for aggregation operations to be scale-invariant; OWA and Choquet aggregations are scale-invariant, while quadratic aggregations aren't.

In Section 4, we explain that OWA and Choquet operations are, in some reasonable sense, the most general ones: namely, crudely speaking, every scale-invariant operation can be composed of linear combinations and min and max operations. We also argue that the selection of linear operations, min, and max as elementary operations is well justified from the computational viewpoint: since they are the fastest possible scale-invariant operations. This justification is presented in Section 5. The mathematical proofs are placed, for reader's convenience, in a special Section 6. The last section contains conclusions.

2 Standard Decision Making Theory: A Brief Reminder

As we have mentioned earlier, to explain our answer to the long-standing puzzle, we need to recall the properties of the utility functions. The needed properties of utility functions are described in this section. Readers who are already well familiar with the standard decision making theory (and with the corresponding properties of utility functions) can skip this section and proceed directly to Section 3.

To be able to describe decisions, we must have a numerical scale for describing preferences. The traditional decision making theory (see, e.g., [5–7]) starts with

an observation that such a scale can be naturally obtained by using probabilities. Specifically, to design this scale, we select two alternatives:

- a very negative alternative A_0 ; e.g., an alternative in which the decision maker loses all his money (and/or loses his health as well), and
- a very positive alternative A_1 ; e.g., an alternative in which the decision maker wins several million dollars.

Based on these two alternatives, we can, for every value $p \in [0, 1]$, consider a randomized alternative $L(p)$ in which we get A_1 with probability p and A_0 with probability $1 - p$.

(It should be mentioned that in the standard decision making theory, randomized alternatives like $L(p)$ are also (somewhat misleadingly) called *lotteries*. This name comes from the fact that a lottery is one of the few real-life examples of randomized outcomes with known probabilities.)

In the two extreme cases $p = 0$ and $p = 1$, the randomized alternative $L(p)$ turns into one of the original alternatives: when $p = 1$, we get the favorable alternative A_1 (with probability 1), and when $p = 0$, we get the unfavorable alternative A_0 . In general, the larger the probability p of the favorable alternative A_1 , the more preferable is the corresponding randomized alternative $L(p)$. Thus, the corresponding randomized alternatives (“lotteries”) $L(p)$ form a continuous 1-D scale ranging from the very negative alternative A_0 to the very positive alternative A_1 .

So, it is reasonable to gauge the preference of an arbitrary alternative A by comparing it to different alternatives $L(p)$ from this scale until we find A 's place on this scale, i.e., the value $p \in [0, 1]$ for which, to this decision maker, the alternative A is equivalent to $L(p)$: $L(p) \sim A$. This value is called the *utility* $u(A)$ of the alternative A in the standard decision making theory.

In our definition, the numerical value of the utility depends on the selection of the alternatives A_0 and A_1 : e.g., A_0 is the alternative whose utility is 0 and A_1 is the alternative whose utility is 1. What if we use a different set of alternatives, e.g., $A'_0 < A_0$ and $A'_1 > A_1$?

Let A be an arbitrary alternative between A_0 and A_1 , and let $u(A)$ be its utility with respect to A_0 and A_1 . In other words, we assume that A is equivalent to the randomized alternative in which:

- we have A_1 with probability $u(A)$, and
- we have A_0 with probability $1 - p$.

In the scale defined by the new alternatives A'_0 and A'_1 , let $u'(A_0)$, $u'(A_1)$, and $u'(A)$ denote the utilities of A_0 , A_1 , and A . This means, in particular:

- that A_0 is equivalent to the randomized alternative in which we get A'_1 with probability $u'(A_0)$ and A'_0 with probability $1 - u'(A_0)$; and
- that A_1 is equivalent to the randomized alternative in which we get A'_1 with probability $u'(A_1)$ and A'_0 with probability $1 - u'(A_1)$.

Thus, the alternative A is equivalent to the compound randomized alternative, in which

- first, we select A_1 or A_0 with probabilities $u(A)$ and $1 - u(A)$, and then
- depending on the first selection, we select A'_1 with probability $u'(A_1)$ or $u'(A_0)$ – and A'_0 with the remaining probability.

As the result of this two-stage process, we get either A'_0 or A'_1 . The probability p of getting A'_1 in this two-stage process can be computed by using the formula of full probability

$$p = u(A) \cdot u'(A_1) + (1 - u(A)) \cdot u'(A_0) = u(A) \cdot (u'(A_1) - u'(A_0)) + u'(A_0).$$

So, the alternative A is equivalent to a randomized alternative in which we get A'_1 with probability p and A'_0 with the remaining probability $1 - p$. By definition of utility, this means that the utility $u'(A)$ of the alternative A in the scale defined by A'_0 and A'_1 is equal to this value p :

$$u'(A) = u(A) \cdot (u'(A_1) - u'(A_0)) + u'(A_0).$$

So, changing the scale means a linear re-scaling of the utility values:

$$u(A) \rightarrow u'(A) = \lambda \cdot u(A) + b$$

for $\lambda = u'(A_1) - u'(A_0) > 0$ and $b = u'(A_0)$.

Vice versa, for every $\lambda > 0$ and b , one can find appropriate events A'_0 and A'_1 for which the re-scaling has exactly these values λ and b . In other words, utility is defined modulo an arbitrary (increasing) linear transformation.

The last important aspect of the standard decision making theory is its description of the results of different actions. Suppose that an action leads to alternatives a_1, \dots, a_m with probabilities p_1, \dots, p_m . We can assume that we have already determined the utility $u_i = u(a_i)$ of each of the alternatives a_1, \dots, a_m . By definition of the utility, this means that for each i , the alternative a_i is equivalent to the randomized alternative $L(u_i)$ in which we get A_1 with probability u_i and A_0 with probability $1 - u_i$. Thus, the results of the action are equivalent to the two-stage process in which, with the probability p_i , we select a randomized alternative $L(u_i)$. In this two-stage process, the results are either A_1 or A_0 . The probability p of getting A_1 in this two-stage process can be computed by using the formula for full probability: $p = p_1 \cdot u_1 + \dots + p_m \cdot u_m$. Thus, the action is equivalent to a randomized alternative in which we get A_1 with probability p and A_0 with the remaining probability $1 - p$. By definition of utility, this means that the utility u of the action in question is equal to

$$u = p_1 \cdot u_1 + \dots + p_m \cdot u_m.$$

In statistics, the right-hand of this formula is known as the *expected value*. Thus, we can conclude that the utility of each action with different possible alternatives is equal to the expected value of the utility.

3 Why Quadratic Aggregation Operations are Less Adequate than OWA and Choquet Operations: An Explanation

To adequately describe the decision maker's preferences, we must be able, given an alternative characterized by n parameters x_1, \dots, x_n , to describe the utility $u(x_1, \dots, x_n)$ of this alternative. To get a perfect description of the user's preference, we must elicit such a utility value for all possible combinations of parameters. As we have mentioned in the Introduction, for practical values n , it is not realistic to elicit that many utility values from a user. So, instead, we elicit the user's preference over each of the parameters x_i , and then aggregate the resulting utility values $u_i(x_i)$ into an approximation for $u(x_1, \dots, x_n)$: $u(x_1, \dots, x_n) \approx f(u_1, \dots, u_n)$, where $u_i \stackrel{\text{def}}{=} u_i(x_i)$.

We have also mentioned that in the first approximation, linear aggregation operations $f(u_1, \dots, u_n) = a_0 + \sum_{i=1}^n w_i \cdot u_i$ work well, but to get a more adequate representation of the user's preferences, we must go beyond linear functions. From the purely mathematical viewpoint, it may seem that quadratic functions $f(u_1, \dots, u_n)$ should provide a reasonable next approximation, but in practice, piece-wise linear aggregation operations such as OWA (or Choquet integral) provide a much more adequate description of expert preferences.

For example, for two parameters, the general OWA combination of two utility values has the form

$$f(u_1, u_2) = w_1 \cdot \min(u_1, u_2) + w_2 \cdot \max(u_1, u_2).$$

Similarly, the general OWA combination of three utility values has the form

$$f(u_1, u_2, u_3) = w_1 \cdot \min(u_1, u_2, u_3) + w_2 \cdot \max(\min(u_1, u_2), \min(u_1, u_3), \min(u_2, u_3)) + w_3 \cdot \max(u_1, u_2, u_3).$$

Let us show that this seemingly mysterious advantage of non-quadratic aggregation operations can be explained based on the main properties of the utility functions.

Indeed, as we have mentioned in Section 2, the utility is defined modulo *two* types of transformations: changing a starting point $u \rightarrow u + b$ and changing a scale $u \rightarrow \lambda \cdot u$ for some $\lambda > 0$. It is therefore reasonable to require that the aggregation operation should not depend on which "unit" (i.e., which extreme event A_1) we use to describe utility. Let us describe this requirement in precise terms.

In the original scale,

- we start with utility values u_1, \dots, u_n ;
- to these values, we apply the aggregation operation $f(u_1, \dots, u_n)$ and get the resulting overall utility $u = f(u_1, \dots, u_n)$.

On the other hand,

- we can express the same utility values in a new scale, as $u'_1 = \lambda \cdot u_1, \dots, u'_n = \lambda \cdot u_n$;
- then, we use the same aggregation function to combine the new utility values; as a result, we get the resulting overall utility $u' = f(u'_1, \dots, u'_n)$.

Substituting the expressions $u'_i = \lambda \cdot u_i$ into this formula, we conclude that $u' = f(\lambda \cdot u_1, \dots, \lambda \cdot u_n)$. We require that the utility

$$u' = f(u'_1, \dots, u'_n) = f(\lambda \cdot u_1, \dots, \lambda \cdot u_n)$$

reflect the same degree of preference as the utility $u = f(u_1, \dots, u_n)$ but in a different scale: $u' = \lambda \cdot u$, i.e.,

$$f(\lambda \cdot u_1, \dots, \lambda \cdot u_n) = \lambda \cdot f(u_1, \dots, u_n).$$

It is worth mentioning that in mathematics, such functions are called *homogeneous* (of first degree). So, we arrive at the conclusion that an adequate aggregation operation should be homogeneous.

This conclusion about the above mysterious fact. On the other hand, one can show that linear aggregation operations and piece-wise linear aggregation operations like OWA are scale-invariant.

Let us start with a linear aggregation operation $f(u_1, \dots, u_n) = w_1 \cdot u_1 + \dots + w_n \cdot u_n$. For this operation, we get

$$\begin{aligned} f(\lambda \cdot u_1, \dots, \lambda \cdot u_n) &= w_1 \cdot (\lambda \cdot u_1) + \dots + w_n \cdot (\lambda \cdot u_n) = \\ &= \lambda \cdot (w_1 \cdot u_1 + \dots + w_n \cdot u_n) = \lambda \cdot f(u_1, \dots, u_n). \end{aligned}$$

Let us now consider the OWA aggregation operation $f(u_1, \dots, u_n) = w_1 \cdot u_{(1)} + \dots + w_n \cdot u_{(n)}$, where $u_{(1)}$ is the largest of n values u_1, \dots, u_n , $u_{(2)}$ is the second largest, etc. If we multiply all the utility values u_i by the same constant $\lambda > 0$, their order does not change. In particular, this means that the same value $u_{(1)}$ which was the largest in the original scale is the largest in the new scale as well. Thus, its numerical value $u'_{(1)}$ can be obtained by re-scaling $u_{(1)}$: $u'_{(1)} = \lambda \cdot u_{(1)}$. Similarly, the same value $u_{(2)}$ which was the second largest in the original scale is the second largest in the new scale as well. Thus, its numerical value $u'_{(2)}$ can be obtained by re-scaling $u_{(2)}$: $u'_{(2)} = \lambda \cdot u_{(2)}$, etc. So, we have $u'_{(i)} = \lambda \cdot u_{(i)}$ for all i . Thus, for the OWA aggregation operation, we have

$$\begin{aligned} f(\lambda \cdot u_1, \dots, \lambda \cdot u_n) &= w_1 \cdot u'_{(1)} + \dots + w_n \cdot u'_{(n)} = w_1 \cdot (\lambda \cdot u_{(1)}) + \dots + w_n \cdot (\lambda \cdot u_{(n)}) = \\ &= \lambda \cdot (w_1 \cdot u_{(1)} + \dots + w_n \cdot u_{(n)}) = \lambda \cdot f(u_1, \dots, u_n). \end{aligned}$$

On the other hand, a generic quadratic operation is not homogeneous. Indeed, a general quadratic operation has the form

$$f(u_1, \dots, u_n) = \sum_{i=1}^n w_i \cdot u_i + \sum_{i=1}^n \sum_{j=1}^n w_{ij} \cdot u_i \cdot u_j.$$

Here,

$$f(\lambda u_1, \dots, \lambda u_n) = \sum_{i=1}^n w_i \cdot (\lambda \cdot u_i) + \sum_{i=1}^n \sum_{j=1}^n w_{ij} \cdot (\lambda \cdot u_i) \cdot (\lambda \cdot u_j) =$$

$$\lambda \cdot \sum_{i=1}^n w_i \cdot u_i + \lambda^2 \cdot \sum_{i=1}^n \sum_{j=1}^n w_{ij} \cdot u_i \cdot u_j.$$

On the other hand,

$$\lambda \cdot f(u_1, \dots, u_n) = \lambda \cdot \sum_{i=1}^n w_i \cdot u_i + \lambda \cdot \sum_{i=1}^n \sum_{j=1}^n w_{ij} \cdot u_i \cdot u_j.$$

The linear terms in the expressions $f(\lambda u_1, \dots, \lambda u_n)$ and $\lambda \cdot f(u_1, \dots, u_n)$ coincide, but the quadratic terms differ: the quadratic term in $f(\lambda u_1, \dots, \lambda u_n)$ differs from the quadratic term in $\lambda \cdot f(u_1, \dots, u_n)$ by a factor of λ . Thus, the only possibility to satisfy the scale-invariance (homogeneity) requirement for all λ is to have these differing quadratic terms equal to 0, i.e., to have $w_{ij} = 0$ – but in this case the aggregation operation is linear. So, quadratic operations are indeed not homogeneous – which explains why they are less adequate in describing user’s preferences than homogeneous operations like OWA or Choquet integral.

4 OWA and Choquet Operations Are, in Some Reasonable Sense, the Most General Ones: A New Result

In the previous section, we explained the empirical fact that in multi-criteria decision making, OWA and Choquet operations lead to more adequate results than seemingly natural quadratic aggregation operations. The explanation is that, due to the known properties of the utility, it is reasonable to require that aggregation operation be scale-invariant (homogeneous); OWA and Choquet operations are scale-invariant but quadratic operations are not.

However, in principle, OWA and Choquet operations are just a few examples of scale-invariant operations, so by itself, the above result does not explain why OWA and Choquet operations are so successful and not any other scale-invariant operation. In this section, we give such an explanation.

This explanation is based on the fact that OWA and Choquet operations are compositions of linear functions, min, and max. In this section, we prove that, crudely speaking, every scale-invariant operation can be composed of linear functions and min and max operations.

Definition 1. A function $f(x_1, \dots, x_n)$ is called homogeneous if for every x_1, \dots, x_n and for every $\lambda > 0$, we have $f(\lambda \cdot x_1, \dots, \lambda \cdot x_n) = \lambda \cdot f(x_1, \dots, x_n)$.

Definition 2. *By a basic function, we mean one of the following functions:*

- a linear function $f(x_1, \dots, x_n) = w_1 \cdot x_1 + \dots + w_n \cdot x_n$;
- a minimum function $f(x_1, \dots, x_n) = \min(x_{i_1}, \dots, x_{i_m})$; and
- a maximum function $f(x_1, \dots, x_n) = \max(x_{i_1}, \dots, x_{i_m})$.

We also say that basic functions are 1-level compositions of basic functions. We say that a function $f(x_1, \dots, x_n)$ is a k -level composition of basic functions if $f(x_1, \dots, x_n) = g(h_1(x_1, \dots, x_n), \dots, h_m(x_1, \dots, x_n))$, where g is a basic function, and the functions $h_1(x_1, \dots, x_n), \dots, h_m(x_1, \dots, x_n)$ are $(k - 1)$ -level compositions of basic functions.

By induction over k , one can easily prove that all compositions of basic functions are homogeneous. For example:

- a linear combination is a basic function;
- an OWA combination of two values is a 2-level composition of basic functions;
- a general OWA operation is a 3-level composition of basic functions.

It turns out that an arbitrary homogeneous function can be approximated by appropriate 3-level compositions.

Definition 3. *Let $k > 0$ be a positive integer. We say that k -level compositions have a universal approximation property for homogeneous functions if for every continuous homogeneous function $f(x_1, \dots, x_n)$, and for every two numbers $\varepsilon > 0$ and $\Delta > 0$, there exists a function $\tilde{f}(x_1, \dots, x_n)$ which is a k -level composition of basic functions and for which $|f(x_1, \dots, x_n) - \tilde{f}(x_1, \dots, x_n)| \leq \varepsilon$ for all x_1, \dots, x_n for which $|x_i| \leq \Delta$ for all i .*

Theorem 1. *3-level compositions have a universal approximation property for homogeneous functions.*

(As we mentioned in Section 1, for readers' convenience, all the proofs are located in the special Proofs section.)

A natural question is: do we need that many levels of composition? What is we only use 1- or 2-level compositions? It turns out that in this case, we will not get the universal approximation property – and thus, the 3 levels of OWA operations is the smallest possible number.

Theorem 2.

- 1-layer computations do not have a universal approximation property for homogeneous functions;
- 2-layer computations do not have a universal approximation property for homogeneous functions.

5 Why Linear Operations, min, and max: A Computational Justification

A natural question is: why should we select linear functions, min, and max as basic functions? One possible answer is that these operations are the fastest to compute, i.e., they require the smallest possible number of computational steps.

Indeed, the fastest computer operations are the ones which are hardware supported, i.e., the ones for which the hardware has been optimized. In modern computers, the hardware supported operations with numbers include elementary arithmetic operations (+, −, ·, /, etc.), and operations min and max.

In the standard (digital) computer (see, e.g., [2])

- addition of two n -bit numbers requires, in the worst case, $2n$ bit operations: n to add corresponding digits, and n to add carries;
- multiplication, in the worst case, means n additions – by each bit of the second factor; so, we need $O(n^2)$ bit operations;
- division is usually performed by trying several multiplications, so it takes even longer than multiplication;
- finally, min and max can be performed bit-wise and thus, require only n bit operations.

Thus, the fastest elementary operations are indeed addition (or, more generally, linear combination), min, and max.

6 Proof of Theorems 1 and 2

1. Before we start proving, let us notice that the values of the functions $\min(x_{i_1}, \dots, x_{i_m})$ and $\max(x_{i_1}, \dots, x_{i_m})$ depend on the order between the values x_1, \dots, x_n . There are $n!$ possible orders, so we can divide the whole n -dimensional space of all possible tuples (x_1, \dots, x_n) into $n!$ zones corresponding to these different orders.

2. In each zone, a basic function is linear:

- a linear function is, of course, linear;
- a minimizing function $\min(x_{i_1}, \dots, x_{i_m})$ is simply equal to the variable x_{i_k} which is the smallest in this zone and is, thus, linear;
- a maximizing function $\max(x_{i_1}, \dots, x_{i_m})$ is simply equal to the variable x_{i_k} which is the largest in this zone and is, thus, also linear.

3. If a function $f(x_1, \dots, x_n)$ can be approximated, with arbitrary accuracy, by functions from a certain class, this means that $f(x_1, \dots, x_n)$ is a limit of functions from this class.

4. Basic functions are linear in each zone; thus, their limits are also linear in each zone. Since some homogeneous functions are non-linear, we can thus conclude that basic functions do not have a universal approximation property for homogeneous functions.

5. Let us now consider 2-level compositions of basic functions, i.e., functions of the type $f(x_1, \dots, x_n) = g(h_1(x_1, \dots, x_n), \dots, h_m(x_1, \dots, x_n))$, where g and h_i are basic functions.

Since there are three types of basic functions, we have three options:

- it is possible that $g(x_1, \dots, x_m)$ is a linear function;

- it is possible that $g(x_1, \dots, x_m)$ is a minimizing function; and
- it is possible that $g(x_1, \dots, x_m)$ is a maximizing function.

Let us consider these three options one by one.

5.1. Let us start with the first option, when $g(x_1, \dots, x_m)$ is a linear function. Since on each zone, each basic function h_i is also linear, the composition $f(x_1, \dots, x_n)$ is linear on each zone.

5.2. If $g(x_1, \dots, x_m)$ is a minimizing function, then on each zone, each h_i is linear and thus, the composition $f(x_1, \dots, x_n)$ is a minimum of linear functions. It is known that minima of linear functions are concave; see, e.g., [8]. So, within this option, the function $f(x_1, \dots, x_n)$ is concave.

5.3. If $g(x_1, \dots, x_m)$ is a maximizing function, then on each zone, each h_i is linear and thus, the composition $f(x_1, \dots, x_n)$ is a maximum of linear functions. It is known that maxima of linear functions are convex; see, e.g., [8]. So, within this option, the function $f(x_1, \dots, x_n)$ is convex.

6. In each zone, 2-level compositions of basic functions are linear, concave, or convex. The class of all functions approximable by such 2-level compositions is the class of limits (closure) of the union of the corresponding three classes: of linear, concave, and convex sets. It is known that the closure of the finite union is the union of the corresponding closures. A limit of linear functions is always linear, a limit of concave functions is concave, and a limit of convex functions is convex. Thus, by using 2-level compositions, we can only approximate linear, concave, or convex functions. Since there exist homogeneous functions which are neither linear nor concave or convex, we can thus conclude that 2-level compositions are not universal approximators for homogeneous functions.

7. To complete the proof, we must show that 3-level compositions are universal approximators for homogeneous functions. There are two ways to prove it.

7.1. First, we can use the known facts about concave and convex functions [8]:

- that every continuous function on a bounded area can be represented as a difference between two convex functions, and
- that every convex function can be represented as a maximum of linear functions – namely, all the linear functions which are smaller than this function.

These facts are true for general (not necessarily homogeneous) functions. For homogeneous functions $f(x_1, \dots, x_n)$, one can easily modify the existing proofs and show:

- that every homogeneous continuous function on a bounded area can be represented as a difference between two convex homogeneous functions, and
- that every homogeneous convex function can be represented as a maximum of homogeneous linear functions – namely, all the homogeneous linear functions which are smaller than this function.

Thus, we can represent the desired function $f(x_1, \dots, x_n)$ as the difference between two convex homogeneous functions $f(x_1, \dots, x_n) = f_1(x_1, \dots, x_n) -$

$f_2(x_1, \dots, x_n)$. Each of these convex functions can be approximated by maxima of linear functions and thus, by 2-level compositions. Substraction $f_1 - f_2$ adds the third level, so $f(x_1, \dots, x_n)$ can indeed be approximated by 3-level compositions.

To prove that a function $f(x_1, \dots, x_n)$ can be represented as a difference between two convex functions, we can, e.g., first approximate it by a homogeneous function which is smooth on a unit sphere $\{(x_1, \dots, x_n) : x_1^2 + \dots + x_n^2 = 1\}$, and then take $f_1(x_1, \dots, x_n) = k \cdot \sqrt{x_1^2 + \dots + x_n^2}$ for a large k . For smooth functions, convexity means that the Hessian matrix – consisting of its second derivatives $\frac{\partial^2 f}{\partial x_i \partial x_j}$ – is positive definite.

For sufficiently large k , the difference

$$f_2(x_1, \dots, x_n) = f_1(x_1, \dots, x_n) - f(x_1, \dots, x_n)$$

is also convex – since its second derivatives matrix is dominated by positive definite terms coming from f_1 . Thus, the difference $f_1 - f_2 = f$ is indeed the desired difference.

7.2. Another, more constructive proof, is, for some $\delta' > 0$, to select a finite δ' -dense set of points $e = (e_1, \dots, e_n)$ on a unit square. For each such point, we build a 2-level composition which coincides with f on the corresponding ray $\{\lambda \cdot (e_1, \dots, e_n) : \lambda > 0\}$. This function can be obtained, e.g., as a minimum of several linear functions which have the right value on this ray but change drastically immediately outside this ray.

For example, let $f_0(x)$ be an arbitrary homogeneous linear function which coincides with $f(x)$ at the point e – and thus, on the whole ray. To construct the corresponding linear functions, we can expand the vector e to an orthonormal basis e, e', e'', \dots , and take linear functions $f_0(x) + k \cdot (e' \cdot x)$ and $f_0(x) - k \cdot (e' \cdot x)$ for all such e' (and for a large $k > 0$). Then, the minimum of all these functions is very small outside the ray.

We then take the maximum of all these minima – a 3-level composition.

The function $f(x_1, \dots, x_n)$ is continuous on a unit sphere and thus, uniformly continuous on it, i.e., for every $\varepsilon > 0$, there is a δ such that δ -close value on the unit sphere lead to ε -close values of f . By selecting appropriate δ' and k (depending on δ), we can show that the resulting maximum is indeed ε -close to f .

The theorem is proven.

7 Conclusions

In multi-criteria decision making, it is necessary to aggregate (combine) utility values corresponding to several criteria (parameters). The simplest way to combine these values is to use linear aggregation. In many practical situations, however, linear aggregation does not fully adequately describe the actual decision making process, so non-linear aggregation is needed.

From the purely mathematical viewpoint, the next natural step after linear functions is the use of quadratic functions. However, in decision making, a different type of non-linearities are usually more adequate than quadratic ones: non-linearities like OWA or Choquet integral that use min and max in addition to linear combinations. In this paper, we explain the empirically observed advantage of such aggregation operations. Specifically, we show that, due to the known properties of the utilities, reasonable operations for aggregating utilities must be scale-invariant. Aggregation operations like OWA or Choquet integral are scale-invariant, while quadratic functions are not. We also prove that operations like OWA and Choquet are general: to be more precise, we prove that compositions of linear functions, min, and max are universal approximators for scale-invariant operators.

Acknowledgments

This work was supported in part by NSF grant HRD-0734825, by Grant 1 T36 GM078000-01 from the National Institutes of Health, by the Japan Advanced Institute of Science and Technology (JAIST) International Joint Research Grant 2006-08, and by the Max Planck Institut für Mathematik.

The authors are very thankful to the anonymous referees for valuable suggestions.

References

1. Ceberio, M., and Modave, F.: An interval-valued, 2-additive Choquet integral for multi-criteria decision making, *Proceedings of the 10th Conf. on Information Processing and Management of Uncertainty in Knowledge-based Systems IPMU'04*, Perugia, Italy, July 2004.
2. Cormen, T.H., Leiserson, C. E., Rivest, R. L., Stein, C.: *Introduction to Algorithms*, MIT Press, Cambridge, Massachusetts, 2001.
3. Feynman, R., Leighton, R., and Sands, M.: *Feynman Lectures on Physics*, Addison Wesley, 2005.
4. Grabisch, M., Murofushi, T., Sugeno, M. eds.: *Fuzzy Measures and Integrals*, Physica-Verlag, Berlin-Heidelberg, 2000.
5. Keeney, R. L., Raiffa, H.: *Decisions with Multiple Objectives*, John Wiley and Sons, New York, 1976.
6. Luce, R. D., Raiffa, H.: *Games and Decisions: Introduction and Critical Survey*, Dover, New York, 1989.
7. Raiffa, H.: *Decision Analysis*, Addison-Wesley, Reading, Massachusetts, 1970.
8. Rockafeller, R. T.: *Convex Analysis*, Princeton University Press, Princeton, New Jersey, 1970.
9. Yager, R. R., Kacprzyk, J., eds.: *The Ordered Weighted Averaging Operators: Theory and Applications*, Kluwer, Norwell, Massachusetts, 1997,